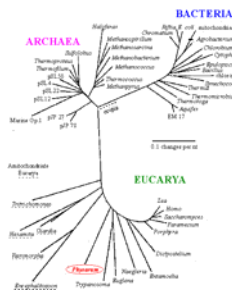


Pairwise Alignment and Database Searching

Anders Gorm Pedersen
Henrik Nielsen
Center for Biological Sequence Analysis

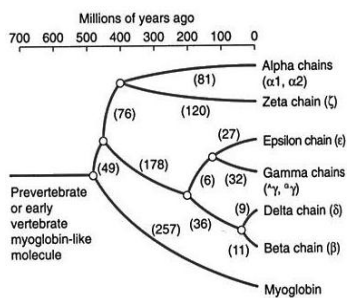
Sequences are related

- Darwin: all organisms are related through descent with modification
- => Sequences are related through descent with modification
- => Similar molecules have similar functions in different organisms



Phylogenetic tree based on ribosomal RNA:
three domains of life

Sequences are related, II



Phylogenetic tree of globin-type proteins found in humans

[illegible]

Pairwise alignment

100.000% identity in 3 aa overlap

SPA
:::
SPA

Percent identity is not a good measure of alignment quality

Pairwise alignments: alignment score

43.2% identity; Global alignment score: 374

	10	20	30	40	50
alpha	V-LSPADKTNVKAANGKVGAGAHGEYGALERMFLSPPTTKTFFPHF-DLS----	HGSA			
beta	VHLTPPEEKSAVTALMGKV--NVDEVGGRALGRLLVVFPWTRFFESFGDLSTPDVAMGNP				
	10	20	30	40	50

	60	70	80	90	100	110
alpha	QVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLKVDFVNFKLLSHCLLVTLAAHL					
beta	KVKARGKKVLGAFSDGLAHLNLIKGTPTATLSBLHCDKLHVDPENFRLGNVLCVLAHNF					
	60	70	80	90	100	110

	120	130	140
alpha	PAEFTPAVHASLDKFLASVSTVLTSKYR		
beta	GEKFTTPVQAAYQKVVGAVANALAHKYH		
	120	130	140

Alignment scores: match vs. mismatch

Simple scoring scheme (too simple in fact...):

Matching amino acids: 5
Mismatch: 0

Scoring example:

K A W S A D V
: : : : :
K D W S A E V
5+0+5+5+5+0+5 = 25

Pairwise alignments: conservative substitutions

43.2% identity; Global alignment score: 374

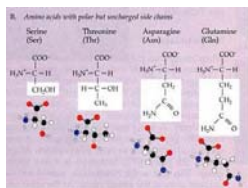
```

alpha  V-LSPADKTNVKAAMGKVGAGAGYGAERMFSPFTTKTYFPHF-DLS----HGSA
      : : : : : : : : : : : : : : : : : : : : : : : : : :
beta   VHLTPDKSAVTALMGKV--NVDEVGGGALGELLVVPWTQRFESPGDLSTPDVAMGNP
      10      20      30      40      50

alpha  QVKGHGKGVADALTNAVAVHVDMPNALSGLDHAHKLKVDPNFKLLSHCLLVTLAAHL
      : : : : : : : : : : : : : : : : : : : : : : : : : :
beta   KVKAHGKVLGAFFSDGLAHLNMLKGTPTATLSLHCDKLHVDPEFNLGNVLCVLAHFF
      60      70      80      90     100     110

alpha  PARETFAVHASLDKFLASVSTVLTSKYR
      : : : : : : : : : : : : : : : : : : : : : : : :
beta   GREFTFPVQAAYQKVVAGVANALAHKYH
      120     130     140
  
```

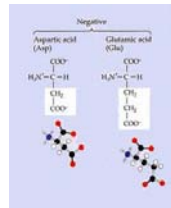
Amino acid properties



Serine (S) and Threonine (T) have similar physicochemical properties

=> Substitution of S/T or E/D occurs relatively often during evolution

=> Substitution of S/T or E/D should result in scores that are only moderately lower than identities



Aspartic acid (D) and Glutamic acid (E) have similar properties

Protein substitution matrices

A	5
R	-2 7
N	-1 -1 7
D	-2 -2 2 8
C	-1 -4 -2 -4 13
Q	-1 1 0 0 -3 7
E	-1 0 0 2 -3 2 6
G	0 -3 0 -1 -3 -2 -3 8
H	-2 0 1 -1 -3 1 0 -2 10
I	-1 -4 -3 -4 -2 -3 -4 -4 5
L	-2 -3 -4 -4 -2 -2 -3 -4 3 2 5
K	-1 3 0 -1 -3 2 1 -2 0 -3 -3 6
M	-1 -2 -2 -4 -2 0 -2 -3 -1 2 3 -2 7
F	-3 -3 -4 -5 -2 -4 -3 -4 -1 0 1 -4 0 8
P	-1 -3 -2 -1 -4 -1 -1 -2 -2 -3 -4 -1 -3 -4 10
S	1 -1 1 0 0 1 0 -1 0 -1 -3 -3 0 -2 -3 -1 5
T	0 -1 0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -2 -1 2 5
W	-3 -3 -4 -5 -5 -1 -3 -3 -3 -3 -2 -3 -1 1 -4 -4 -3 15
Y	-2 -1 -2 -3 -3 -1 -2 -3 2 -1 -1 -2 0 4 -3 -2 -2 2 8
V	0 -3 -3 -4 -1 -3 -3 -4 -4 4 1 -3 1 -1 -3 -2 0 -3 -1 5

BLOSUM50 matrix:

- Positive scores on diagonal (identities)
- Similar residues get higher (positive) scores
- Dissimilar residues get smaller (negative) scores

Pairwise alignments: insertions/deletions

43.2% identity; Global alignment score: 374

```
alpha  V-LSPADKTNVKAAMGKVGAGAHGEYGAEALERMFLSFPTTKTYFPHF-DLS---HGSA
      10 20 30 40 50
beta   VHLTPEEKSAVTALMGKV--NVDEVGGEALGRLLVVYPWTQFFESFGDLSTPDAVMGNP
      10 20 30 40 50

alpha  QVKGHGKGVADALTNVAHVDDMPNALSGLSDLAHKLKVDQPNFKLLSHCLLVTLAAHL
      60 70 80 90 100 110
beta   KVKAGKGVLGAFSDGLAHLNMLKGTFTLSELMCDKLVDPENFRLGNVLCVLAHMF
      60 70 80 90 100 110

alpha  PARFTPAVHASLDKFLASVSTVLTSKYR
      120 130 140
beta   GREFTFPVQAAYGKVVAGVANAIAHKYH
      120 130 140
```

Alignment scores: insertions/deletions

```
K L A A S V I L S D A L
K L A A - - - S D A L
```

$-10 + 3 \times (-1) = -13$

Affine gap penalties:
Multiple insertions/deletions may be one evolutionary event =>
Separate penalties for gap opening and gap elongation

Handout

Compute 4 alignment scores: two different alignments using two different alignment matrices (and the same gap penalty system)

- Score 1: Alignment 1 + BLOSUM-50 matrix + gaps
- Score 2: Alignment 1 + BLOSUM-Trp matrix + gaps
- Score 3: Alignment 2 + BLOSUM-50 matrix + gaps
- Score 4: Alignment 2 + BLOSUM-Trp matrix + gaps

Note: fake matrix constructed for pedagogic purposes.

Estimation of an empirical matrix

	60	70	80	90	100	110
alpha	QVKGHGKKVADALTNAVAHVDDMPNALSASDLHAHKLSDVDPVNFKLLSHCLLVTLAAHL					

beta	KVKAHGKKVLGAFSDGLAHDLNLTGTATLSLHCDKLEVDPENFLLGNVLVCVLAHHF					
	60	70	80	90	100	110

- Start from given alignments of closely related proteins
- Count the aligned amino acid pairs (e.g., A aligned with A makes up 1.5% of all pairs. A aligned with C makes up 0.01% of all pairs, etc.)
- Expected pair frequencies are computed from single amino acid frequencies. (e.g, $f_{A,C} = f_A \times f_C = 7\% \times 3\% = 0.21\%$).
- For each amino acid pair the substitution scores are essentially computed as:

$$\log \frac{\text{Pair-freq(obs)}}{\text{Pair-freq(expected)}} \quad S_{A,C} = \log \frac{0.01\%}{0.21\%} = -1.3$$

- To obtain the PAM1 matrix, normalize pair frequencies to 1% difference before applying the logarithm
- To obtain the PAM50 matrix, extrapolate the PAM1 matrix via matrix multiplication

Estimation of the BLOSUM 50 matrix

- Use the BLOCKS database (ungapped alignments of especially conserved regions of multiple alignments)
- For each alignment in the BLOCKS database the sequences are grouped into clusters with at least 50% identical residues (for BLOSUM 50)
- All pairs of sequences are compared *between* clusters, and the observed pair frequencies are noted
- Substitution scores are calculated as before

ID	FIBRONECTIN_3; BLOCK
COG9_CANFA	QNSAGSPCVFFFIPLGKQSTCTTREGSDGMLCATT
COG9_RABIT	QNSAGSPCHFFFTFEGSEYACTTDGSDGMACSTT
PA11_HUMAN	LTVYDSCHFFFPQVHSLGHECTRESDPQPMWCATT
HSPA_HUMAN	LTFDGSCHFFFPYDSMLACTESGASGSENCATSE
MANR_HUMAN	QNSAGATCAFFFKFENWTADCTSGSDGMWCUTT
MPR1_MOUSE	RTDGDSPCVFFFIYKQSYDSGLVSGSALMCSTAN
SP1_PIG	ALTDSGKCVFFFIYKHLFDGLSDTYMCVCVTV
SPF1_BOVIN	ELPDESECVFFFPYKSKFQCTVBSGLFWCLSDAD
SPF3_BOVIN	ARTENKSCVFFFIYKSKYFQCTLBSGLFWCLSDAD
SPF4_BOVIN	AUTPDACAFPPFYKSKYFQCTLBSGLFWCLSDAD
SP1_BOISE	AATDVAKCAFPPFYKQYQCTDGLFLRWMCVTV
COG2_CHICK	QNSGSPCVFFFIPLONKTSCTSGSDGMWCATT
COG2_HUMAN	QNSGSPCVFFFIPLONKTSCTSGSDGMWCATT
COG2_MOUSE	QNSGSPCVFFFIPLONKTSCTSGSDGMWCATT
COG2_RABIT	QNSGSPCVFFFIPLONKTSCTSGSDGMWCATT
COG2_RAT	QNSGSPCVFFFIPLONKTSCTSGSDGMWCATT
COG9_BOVIN	QNSAGSPCVFFFIPLGKQSTCTTREGSDGMLCATT
COG9_HUMAN	QNSAGSPCVFFFIPLGKQSTCTTREGSDGMLCATT
COG9_MOUSE	QNSAGSPCVFFFIPLGKQSTCTTREGSDGMLCATT
COG9_RAT	QNSAGSPCVFFFIPLGKQSTCTTREGSDGMLCATT
FINC_BOVIN	QNSGALCHFFFIYKSKYFQCTLBSGLFWCLSDAD
FINC_HUMAN	QNSGALCHFFFIYKSKYFQCTLBSGLFWCLSDAD
FINC_RAT	QNSGALCHFFFIYKSKYFQCTLBSGLFWCLSDAD
MPR1_BOVIN	RTDGDSPCVFFFIYKQSYDSGLVSGSALMCSTAN
MPR1_HUMAN	RTDGDSPCVFFFIYKQSYDSGLVSGSALMCSTAN
PA12_BOVIN	QNSAGTCHFFFIPLGKQSTCTTREGSDGMLCATT
PA12_RABIT	QNSAGTCHFFFIPLGKQSTCTTREGSDGMLCATT

Pairwise alignment

Optimal alignment:

alignment having the highest possible score given a substitution matrix and a set of gap penalties

So:

best alignment can be found by exhaustively searching all possible alignments, scoring each of them and choosing the one with the highest score?

The problem:
How many possible alignments are there?

ACG AC-G --ACG -A-CG
ACG ACG- AC-G- A-CG-

-ACG AC-G --ACG ...
ACG- A-CG A-CG-

-ACG AC-G --ACG
AC-G -ACG AC--G

-ACG ACG- --ACG
A-CG AC-G A-C-G

A-CG ACG- --ACG
ACG- A-CG A--CG

A-CG ACG- -A-CG
AC-G -ACG ACG--

A-CG --ACG -A-CG
-ACG ACG-- AC-G-

Pairwise alignment: the problem

The number of possible pairwise alignments increases explosively with the length of the sequences:

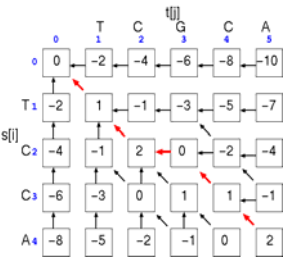
$$f(n_1, n_2) = \sum_{i=0}^{n_1} \binom{n_1}{i} \binom{n_2 + i}{n_1}$$

Two protein sequences of length 100 amino acids can be aligned in approximately 10^{60} different ways

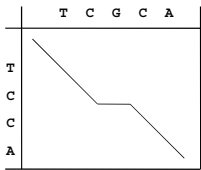
Time needed to test all possibilities is same order of magnitude as the entire lifetime of the universe.

Pairwise alignment: the solution

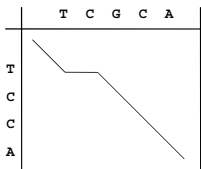
"Dynamic programming"
(the Needleman-VWunsch algorithm)



Alignment depicted as path in matrix

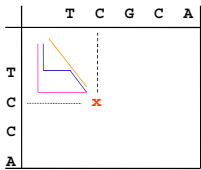


TCGCA
TC-CA



TCGCA
T-CCA

Alignment depicted as path in matrix



Meaning of point in matrix:
all residues up to this point
have been aligned (but there
are many different possible
paths).



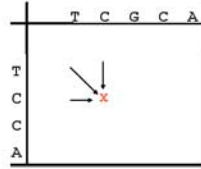
Position labeled "x": TC aligned with TC

--TC
TC--

-TC
T-C

TC
TC

Dynamic programming: computation of scores



Any given point in matrix can only be
reached from three possible previous
positions (you cannot "align
backwards").

=> Best scoring alignment ending in
any given point in the matrix can be
found by choosing the highest
scoring of the three possibilities.

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C					
C					
A					

Any given point in matrix can only be reached from three possible positions (you cannot "align backwards").

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \end{cases}$$

Dynamic programming: computation of scores

	T	C	G	C	A
T					
C					
C					
A					

Any given point in matrix can only be reached from three possible positions (you cannot "align backwards").

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \end{cases}$$

Dynamic programming: computation of scores

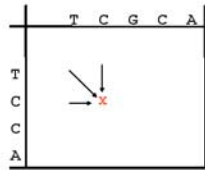
	T	C	G	C	A
T					
C					
C					
A					

Any given point in matrix can only be reached from three possible positions (you cannot "align backwards").

=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \\ \text{score}(x-1,y) - \text{gap-penalty} \end{cases}$$

Dynamic programming: computation of scores



Any given point in matrix can only be reached from three possible positions (you cannot "align backwards").

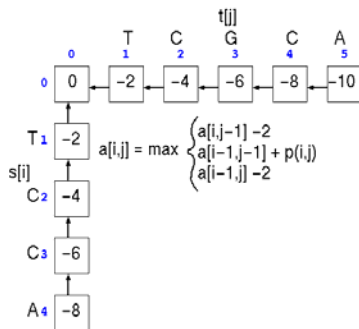
=> Best scoring alignment ending in any given point in the matrix can be found by choosing the highest scoring of the three possibilities.

Each new score is found by choosing the maximum of three possibilities. For each square in matrix: keep track of where best score came from.

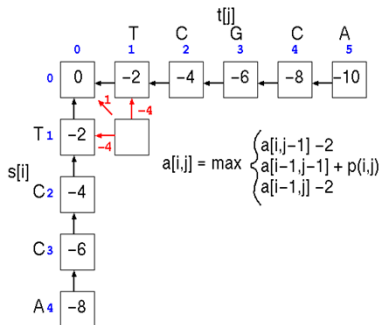
Fill in scores one row at a time, starting in upper left corner of matrix, ending in lower right corner.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \\ \text{score}(x-1,y) - \text{gap-penalty} \end{cases}$$

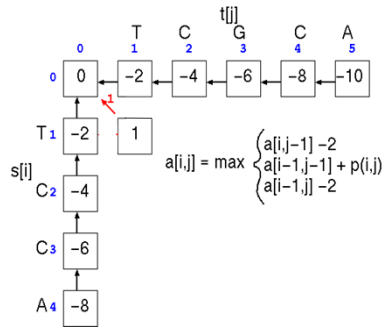
Dynamic programming: example



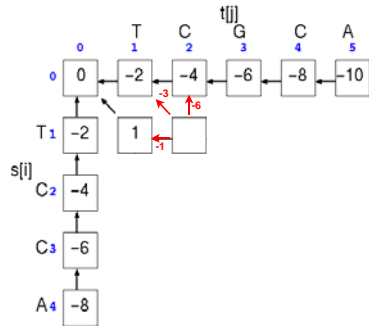
Dynamic programming: example



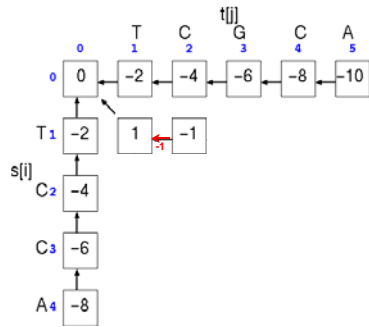
Dynamic programming: example



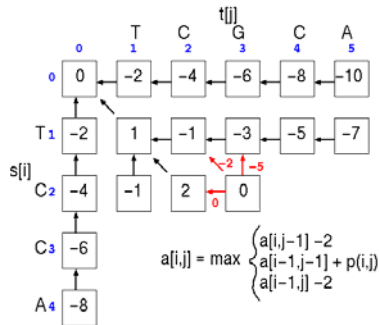
Dynamic programming: example



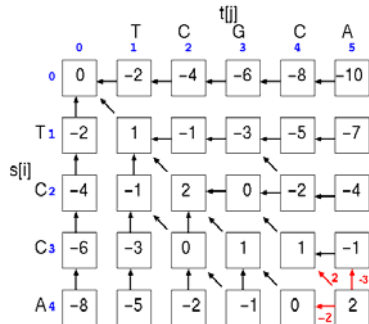
Dynamic programming: example



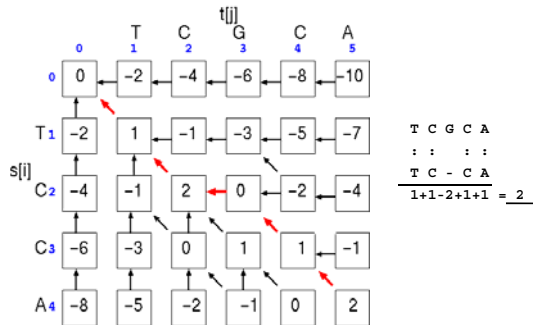
Dynamic programming: example



Dynamic programming: example



Dynamic programming: example

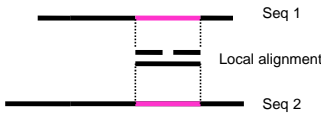


Global versus local alignments

Global alignment: align full length of both sequences.
(The "Needleman-Wunsch" algorithm).

 Global alignment

Local alignment: find best partial alignment of two sequences
(the "Smith-Waterman" algorithm).

 Seq 1
Local alignment
Seq 2

Local alignment overview

• The recursive formula is changed by adding a fourth possibility: zero. This means local alignment scores are never negative.

$$\text{score}(x,y) = \max \begin{cases} \text{score}(x,y-1) - \text{gap-penalty} \\ \text{score}(x-1,y-1) + \text{substitution-score}(x,y) \\ \text{score}(x-1,y) - \text{gap-penalty} \\ 0 \end{cases}$$

- Trace-back is started at the highest value rather than in lower right corner
- Trace-back is stopped as soon as a zero is encountered

Local alignment: example

		H	E	A	G	A	W	G	H	E	E
P	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	0	0	0
W	0	0	0	0	2	0	20	12	4	0	0
H	0	10	2	0	0	0	12	18	22	14	6
E	0	2	16	8	0	0	4	10	18	28	20
A	0	0	8	21	13	5	0	4	10	20	27
E	0	0	6	13	18	12	4	0	4	16	26

AWGHE
AW-HE

Substitution matrices and sequence similarity

Substitution matrices come as series of matrices calculated for different degrees of sequence similarity (different evolutionary distances).

"Hard" matrices	"Soft" matrices
Designed for very similar sequences	Designed for less similar sequences
High numbers in the BLOSUM series (e.g., BLOSUM90)	Low numbers in the BLOSUM series (e.g., BLOSUM30)
Low numbers in the PAM series (e.g. PAM30)	High numbers in the PAM series (e.g. PAM250)
Severe mismatch penalties	Less severe mismatch penalties
Yield short alignments with high %identity	Yield longer alignments with lower %identity

Alignments: things to keep in mind

"Optimal alignment" means "having the highest possible score, given substitution matrix and set of gap penalties".

This is NOT necessarily the biologically most meaningful alignment.

Specifically, the underlying assumptions are often wrong: substitutions are not equally frequent at all positions, affine gap penalties do not model insertion/deletion well, etc.

Pairwise alignment programs always produce an alignment - even when it does not make sense to align sequences.
